

UNITED STATES PATENT APPLICATION

**REMOTE TRANSLATION MECHANISM FOR A MULTINODE
SYSTEM**

INVENTORS:

Kitrick Sheets

Citizenship: USA Residence: **Morrisville, NC**

Post Office Address: **107 Monument View Lane, Morrisville, NC 27560**

Andrew B. Hastings

Citizenship: USA Residence: **Mendota Heights, MN**

Post Office Address: **1340 Mendota Heights Road, Mendota Heights, MN 55120**

1 **REMOTE TRANSLATION MECHANISM FOR A MULTINODE**
2 **SYSTEM**

3
4 **Related Invention**

5 The present invention is related to U.S. Patent Application Serial Number
6 10/235,898, entitled "Remote Translation Mechanism for a Multi-Node System"
7 filed on September 4, 2002, which is incorporated in its entirety herein by reference.

8
9 **Field of the Invention**

10 The present invention relates generally to the field of computer memory
11 systems, and more particularly to a remote translation mechanism for a multinode
12 system.

13
14 **Background of the Invention**

15 Multiprocessor computer systems include a number of processing nodes
16 connected together by an interconnection network. Typically, each processing node
17 includes one or more processors, a local memory, and an interface circuit
18 connecting the node to the interconnection network. The interconnection network is
19 used for transmitting packets of information between processing nodes.
20 Distributed, shared-memory multiprocessor systems include a number of processing
21 nodes that share a distributed memory element. By increasing the number of
22 processing nodes, or the number of processors within each node, such systems can
23 often be scaled to handle increased demand. In such a system, each processor is
24 able to access local memory, or memory of other (remote) processing nodes.
25 Typically, a virtual address is used for all memory accesses within a distributed,
26 shared-memory multiprocessor system, and is translated into a physical address in
27 the requesting node's translation look-aside buffer (TLB). Thus, the requesting
28 node's TLB will need to contain address translation information for all the memory
29 that the node is able to access (local or remote). This amount of address translation

1 information can be substantial, and can result in much duplication of translation
2 information throughout the multiprocessor system (e.g., if the same page of memory
3 is accessed by 64 different nodes, the TLB used by each node will need to contain
4 an entry for that page). This type of system does not scale efficiently to very large
5 memories.

6 Therefore, there is a need for an address translation mechanism in a multi-
7 processor system that addresses these and other shortcomings.

8

9 **Summary of the Invention**

10 To address these and other needs, various embodiments of the present
11 invention are provided. One embodiment of the invention provides a method for
12 initializing shared memory in a multinode system. The method includes building a
13 local address space in each of a plurality of nodes and exporting the local address
14 space from each of the plurality of nodes to a Remote Translation Table (RTT) in
15 each of the plurality of nodes.

16 Another embodiment of the invention provides a method for remotely
17 translating a virtual memory address into a physical memory address in a multi-node
18 system. The method includes providing the virtual memory address at a source
19 node, determining that the virtual memory address is to be sent to a remote node,
20 sending the virtual memory address to the remote node, and translating the virtual
21 memory address on the remote node into a physical memory address using a
22 remote-translation table (RTT). The RTT contains translation information for an
23 entire virtual memory address space associated with the remote node.

24 A further embodiment of the invention provides a method for translating a
25 virtual memory address in a multi-node system. The method includes providing a
26 virtual memory address on a local node by using a virtual address of a load or a
27 store instruction, identifying a virtual node associated with the virtual memory
28 address, and determining if the virtual node corresponds to the local node. If the
29 virtual node corresponds to the local node, then the method includes translating the
30 virtual memory address into a local physical memory address on the local node. If,

1 instead, the virtual node corresponds to a remote node, then the method includes
2 sending the virtual memory address to the remote node, and translating the virtual
3 memory address into a physical memory address on the remote node.

4 These and other embodiments will be described in the detailed description below.

5

6

Brief Description of the Drawings

7

FIG. 1 illustrates a block diagram of a node that includes four multi-streaming
8 processors, according to one embodiment of the present invention.

8

9

FIG. 2A illustrates a format for a physical memory address, according to one
10 embodiment of the present invention.

10

11

FIG. 2B illustrates a more detailed map of the physical memory address format
12 shown in FIG. 2A, according to one embodiment of the present invention.

12

13

FIG. 3 illustrates a format for a virtual memory address, according to one
14 embodiment of the present invention.

14

15

FIG. 4 illustrates a flow diagram for analyzing a VNode field in the virtual memory
16 address, according to one embodiment of the present invention.

16

17

FIG. 5 illustrates a detailed flow diagram for implementing remote translation of a
18 virtual memory address, according to one embodiment of the present
19 invention.

19

20

FIG. 6A illustrates a block diagram of a remote-translation table (RTT) resident on
21 an M chip, according to one embodiment of the present invention.

21

22

FIG. 6B illustrates a more detailed diagram of the RTT shown in FIG. 6A,
23 according to one embodiment of the present invention.

23

24

FIG. 6C illustrates a more detailed view of a portion of the RTT shown in FIG. 6B,
25 according to one embodiment of the present invention.

25

26

FIG. 6D illustrates an even further detailed view of a portion of the RTT shown in
27 FIG. 6C, according to one embodiment of the present invention.

27

28

FIG. 7 provides a functional view of a model for routing memory access requests
29 according to an embodiment of present invention.

29

1 FIG. 8 provides a method for construction of a distributed address space according
2 to an embodiment of the present invention.
3 FIG. 9 provides a flow diagram of RTT aggregation according to an embodiment of
4 the present invention.
5 FIG. 10 shows a flow diagram of a method according to an embodiment of the
6 present invention.
7 FIG. 11 shows a perspective view of a system according to an embodiment of the
8 present invention.
9

10 **Detailed Description**

11 In the following detailed description of the preferred embodiments, reference
12 is made to the accompanying drawings that form a part hereof, and in which are
13 shown by way of illustration specific embodiments in which the invention may be
14 practiced. It is understood that other embodiments may be utilized and structural
15 changes may be made without departing from the scope of the present invention.

16 The leading digit(s) of reference numbers appearing in the Figures generally
17 corresponds to the Figure number in which that component is first introduced, such
18 that the same reference number is used throughout to refer to an identical
19 component which appears in multiple Figures. The same reference number or label
20 may refer to signals and connections, and the actual meaning will be clear from its
21 use in the context of the description.

22 Various embodiments of the present invention provide a virtual-to-physical
23 address translation mechanism for a shared-memory multiprocessor that scales
24 efficiently to large numbers of processors. This mechanism supports a single virtual
25 address format (e.g., using load or store instructions), and detects whether a
26 reference for the instruction is to the local node or a remote node. If to a local node,
27 or if remote translation is not enabled, the virtual to physical address translation is
28 performed in the local translation look-aside buffer (TLB), producing a physical
29 address that includes both the physical node number and the physical memory offset
30 within that node. If remote translation is enabled, however, and the virtual address

1 is for a remote node (as determined by comparing the virtual node field of the
 2 virtual address with the value representing the local virtual node number), then a
 3 remote address translation mechanism is used, as follows. A physical node number
 4 is determined by adding the virtual node field of the virtual address to a physical
 5 base node. The virtual node number is also checked against a limit value, to ensure
 6 that the request is within allowable bounds. The remainder of the virtual address
 7 forms a virtual offset, which is sent with the memory request to the destination
 8 physical node. A "global address space identifier" (GASID) is also looked up for the
 9 local processor and sent with the request. The GASID and the upper portion of the
 10 virtual address are used to index into a remote translation table (RTT) at the
 11 destination node, to produce a physical page number at the remote node. The RTT
 12 is sized to cover the entire virtual address space at a single node. The use of the
 13 GASID allows multiple applications, with overlapping virtual address ranges, to
 14 share memory on the same node while all using the remote translation mechanism.
 15 Essentially, the GASID forms a unique extension to the virtual address offset for
 16 each application.

17 The address translation mechanism of these embodiments scales to large
 18 system sizes, because each node keeps track of virtual to physical page mappings
 19 for its node only. The TLB is used for references to the local node by the local
 20 processor, and the RTT at a node is used for incoming references to the local node
 21 from remote nodes. A single virtual address format and access mechanism are used
 22 for both local and remote memory references. The use of remote translation is thus
 23 functionally transparent. The RTT keeps a full map of the remote virtual address
 24 space, and each node is able to manage its virtual-to-physical address mapping
 25 independently.

26 FIG. 1 illustrates a specific hardware environment in which various
 27 embodiments of the present invention may be practiced. It is to be noted that FIG. 1
 28 illustrates only one example of a hardware environment, and other environments
 29 (for other embodiments) may also be used.

1 FIG. 1 illustrates a block diagram of a node that includes four multi-
2 streaming processors (MSP's), according to one embodiment. In this embodiment,
3 node 100 includes each MSP 102 in a four MSP system. Node 100 is contained on
4 a single printed circuit board. The sixteen M chips on node 100 contain memory
5 controllers, network interfaces and cache coherence directories with their associated
6 protocol engines. The memory system is sliced across the 16 M chips, round robin
7 by 32-byte cache lines. Each M chip supports one slice. Bits 5 and 6 of the
8 physical address determine the E chip within a processor, and bits 7 and 8 further
9 specify one of four M chips connected to each E chip.

10 Each M chip resides in one of sixteen independent address slices of the
11 machine, and the interconnection network provides connectivity only between
12 corresponding M chips on different nodes. All activity (cache, memory, network)
13 relating to a line of memory stays within the corresponding slice. Each M chip
14 controls a separate sector of a slice. Slices expand (get more memory in each) as
15 nodes are added so the number of sectors in each slice is equal to the number of
16 nodes in a system.

17 Total peak local memory bandwidth for one node is 204.8 GB/s, or 51.2
18 GB/s per MSP. As each MSP 102 needs a maximum bandwidth of about 45 GB/s,
19 there is bandwidth to support network traffic and I/O without greatly impacting
20 computational performance. Each M chip contains two network ports, each 1.6
21 GB/s peak per direction.

22 Node 100 also contains two I chip I/O controller ASIC's. These connect to
23 the M chips and provide four I/O ports of 1.2 GB/s bandwidth, full duplex, off node
24 100. Each I chip contains two ports, 400 MB/s full duplex connections to 8 of the
25 local M chips (one I chip connects to the even M chips and the other connects to the
26 odd M chips), and a 1.6 GB/s full duplex connection to the other I chip. The total
27 I/O bandwidth per module is thus 4.8 GB/s full duplex.

28 The memory on node 100 is distributed across the set of 16 M chips. Each
29 M chip directly controls the resources contained on two daughter boards so that
30 there are thirty two daughter boards on node 100. The memory chips in the

1 daughter boards are Direct Rambus DRAM. These chips have 16 internal banks and
2 have 18 data pins that each run, with a 400 MHz clock, at an 800 Mbaud rate. Each
3 chip then has a 1.6 GB/s read/write data rate. Being 18 bits wide, additional parts to
4 support ECC are not needed. Daughter cards contain 16 chips organized into 4
5 memory channels of 4 chips each. Each memory channel is independent. Channels
6 have a peak data bandwidth of 1.6 GB/s, so that the card supports a bandwidth of
7 6.4 GB/s. With 16 banks in a memory chip, a channel has 64 memory banks.
8 Daughter cards with 64 Mbit, 128 Mbit, 256 Mbit or 512 Mbit chips are supported.
9 The design also accommodates chip densities of 1 Gbit if and when they become
10 available, assuming they fit in the design envelope (size, power, etc.). As the
11 memory parts are 18 bits wide instead of 16 in order to support ECC, the chip's bit
12 densities are actually 72, 144, 288, 576 and 1152 Mbits.

13 FIG. 2A illustrates a format for a physical memory address, according to one
14 embodiment. In this embodiment, a 46-bit (64 TBytes) physical memory address is
15 supported. The node size for this embodiment is a board containing four MSP's and
16 16 M chips. Physical memory address format 200 contains bits 47..0. Bits 35..0
17 represent an offset (into memory). Bits 45..36 represent the node. Bits 47..46
18 represent the physical address space. The physical memory format allows for up to
19 1024 nodes (4096 MSP's) and 64 GBytes of physical memory per node. Physical
20 pages are allocated on a per-node basis. That is, any given physical page is
21 distributed uniformly across the 16 sectors (the memory controlled by a given M
22 chip) of a single node. This embodiment provides three parallel, physical address
23 spaces, which are selected by two extra bits at the top of the physical address.

24 FIG. 2B illustrates a more detailed map of the physical memory address
25 format shown in FIG. 2A, in one embodiment. The mapping of a physical address
26 to a destination location is dependent on the hardware implementation (as opposed
27 to being specified in the user-level architecture). Address mapping must be done so
28 that parallelism can be easily exploited by having the map such that multiple
29 transactions can be requested and satisfied simultaneously with minimum hardware
30 complexity. Bits 4..0 represent the byte in the line. Bits 6..5 represent the quadrant

1 (E chip). Bits 8..5 collectively represent the slice/section (M chip). Bits 11..9
 2 represent the memory channel. Bits 13..12 represent the memory chip for the
 3 memory channel, and bits 17..14 represent the bank for the memory chip. Bits
 4 35..18 represent the memory chip address, and bits 45..36 represent the node
 5 number (in the system). Bits 47..46 represent the address space. Memory size
 6 options and configuration changes (including memory upgrades) can modify this
 7 map. The map supports memory chips up to 1 Gbit density. There are three address
 8 spaces: coherent main memory, memory-mapped register space, and I/O device
 9 space. Coherent main memory may be cached.

10 FIG. 3 illustrates a format for a virtual memory address, according to one
 11 embodiment. In this embodiment, virtual memory address format 300 contains a
 12 64-bit virtual address space. Bits 37..0 represent a virtual offset into virtual memory
 13 space, wherein potential page boundaries range from 64 KB to 4 GB. Bits 47..38
 14 represent the VNode (i.e., virtual node). This is used by the hardware when
 15 performing remote address translation. Bits 61..48 must be set to zero in this
 16 implementation. Bits 63..62 specify the memory region, which determines the type
 17 of address translation used in kernel mode. The virtual address space can be
 18 considered a flat virtual address space for uniprocessor, or symmetric
 19 multiprocessing applications. As stated, this embodiment supports eight page sizes
 20 ranging from 64 KB to 4 GB. Thus, the page boundary can vary, from between bits
 21 15 and 16, to between bits 31 and 32.

22 In various embodiments of the invention, virtual addresses used for
 23 instruction fetches and data references are first translated into physical addresses
 24 before memory is accessed. These embodiments support two forms of address
 25 translation: source translation, and remote translation. The first form of address
 26 translation is source translation, in which a virtual address is fully translated by a
 27 Translation Look-aside Buffer (TLB) on a local P chip to a physical address on an
 28 arbitrary node. The second form of address translation is remote translation, in
 29 which the physical node number is determined by a simple translation of the virtual
 30 address VNode field, and the remaining virtual address VOffset field is sent to the

1 remote node to be translated into a physical address offset via a Remote-Translation
 2 Table (RTT). The type of address translation performed is based upon values in a
 3 configuration control register and the virtual address itself. Remote translation is
 4 performed if all of the following three conditions are true: (1) Remote translation is
 5 enabled (e.g., a flag contained in the configuration control register is set); (2) The
 6 virtual address is to the useg region (Bits 63..62 = 00 in the virtual address); and (3)
 7 The virtual address references a remote node (Bits 47..38 in the virtual address are
 8 not equal to a local node value contained in the configuration control register). If
 9 any of the above conditions are false, then source translation is performed. Remote
 10 translation can be enabled/disabled on a per-processor basis.

11 FIG. 4 illustrates a flow diagram for analyzing a VNode field in the virtual
 12 memory address, according to one embodiment of the present invention. Flow
 13 diagram 400 includes blocks 402, 406, and 408, and also includes checkpoint 404.
 14 Flow diagram 400 illustrates one way in which a virtual memory address can be
 15 translated into a physical memory address (in either local or remote memory space).
 16 Block 402 includes identifying the virtual node from a virtual address. In one
 17 implementation, a local node can identify the virtual node by looking at the VNode
 18 field of the virtual address. Checkpoint 404 determines if the virtual node is the
 19 same as, or equal to, the local node. If so, flow diagram 400 continues to block 406,
 20 wherein the virtual address is translated into a physical address locally using a
 21 Translation Look-Aside Buffer (TLB). The local node is then able to address local
 22 physical memory space. If the virtual node is not the same as the local node, then
 23 flow diagram 400 continues to block 408, wherein the virtual address is translated
 24 into a physical address remotely (on a remote node) using a Remote-Translation
 25 Table (RTT). In this fashion, the local node is effectively able to address remote
 26 memory space of the remote node.

27 FIG. 5 illustrates a detailed flow diagram 500 for implementing remote
 28 translation of a virtual memory address, according to one embodiment of the present
 29 invention. When remote translation is enabled, the hardware treats bits 47..38 of the
 30 Virtual Address (i.e., VA47..38, for the VNode field) as a virtual node number. As

1 described above, all use virtual addresses with a VNode value not matching the
2 local virtual node number are translated remotely. Additionally, the address is
3 checked to make sure it does not exceed the user's partition defined by the
4 NodeLimit field in the TLBcontrol register. If $VA_{47..38} > \text{NodeLimit}$, then an
5 Address Error exception occurs.

6 The physical node number for the address is computed by adding $VA_{47..38}$
7 to the BaseNode value from the TLBcontrol register. (In this instance, the
8 BaseNode is a reference mechanism by which the physical node number can be
9 computed.) Overflow on this 10-bit addition is ignored; the OS must never create a
10 partition (via the BaseNode and NodeLimit values) that exceeds the number of
11 nodes in the machine. The virtual address Offset field ($VA_{37..0}$) is sent to the
12 resulting physical node as an Aremote virtual address@ (RVA) to complete the
13 translation. (The RVA could also be referred to as a remote virtual memory
14 address.) The cache allocation is forced to non-allocate and the reference is not
15 cached (Get/Put semantics).

16 The value of the BaseNode is unique to each node. This creates a unique
17 physical node mapping when adding the BaseNode to the VNode field. Therefore,
18 in one implementation, various nodes can use common, contiguous VNodes
19 (starting at 0, for example) to effectively reference different physical nodes (that are
20 used for routing). Table 1 below illustrates an example of physical (destination)
21 node mapping for three different source nodes A, B, and C.

Source Node	VNode	Physical Node
A (Base Node = 100)	0	100
A	1	101
A	2	102
B (Base Node = 200)	0	200
B	1	201
B	2	202
C (Base Node = 300)	0	300
C	1	301
C	2	302

Table 1 – Physical Node Calculation

1 In another embodiment, a look-up table is used to determine the physical
2 node. In this embodiment, the BaseNode calculation is not required.

3 RVA requests bypass the Ecache (in the E chips), since they can never be
4 cached. The M chips contain a set of four, 2-bit Global Address Space ID (GASID)
5 registers, one for each of the local MSP's. When the local M chip sends a packet out
6 the network with an RVA, it includes the value of the two bit GASID for the
7 originating MSP. This is used to qualify the remote translation of the RVA at the

1 destination M chip. Thus, the 2-bit GASID, and the RVA, are routed through the
2 interconnection network. Bits 8..0 of the virtual address are not routed as such,
3 because bits 8..5 are used to select the memory/network slice, and bits 4..0 are used
4 to generate the cache line word mask. The cache line mask is unneeded for the
5 remote translation mechanism, and the slice information is also unneeded, since the
6 remote transfer operations have an implied slice (from an M chip on one node to a
7 corresponding M chip on the remote node).

8 At the remote M chip, remote virtual addresses go through a translation to a
9 pure physical address. This translation takes place before presenting the packet to
10 the directory protocol engine. Remote translation takes place with a granularity of
11 16 MB. The two GASID bits, and bits 37..24 of the RVA, are used to index into a
12 64K-entry Remote-Translation Table (RTT). Each entry of this table contains a
13 valid bit, a write-enable bit, and a 12-bit value representing PA35..24 (the 16 MB
14 physical page frame). These bits are appended to the lower bits of the RVA to form
15 a physical memory address at the remote node. The valid bit is used for status of the
16 translation. The valid bit indicates whether the RTT was able to translate the virtual
17 memory address into a valid physical memory address space on the remote node.

18 The write-enable bit, or flag, indicates whether a write is permitted to a
19 region referenced by the virtual memory address. A write to this region will only be
20 allowed if the write-enable bit is set.

21 In one embodiment, the RVA is formed from one or more portions of the
22 virtual address having the VNode field. In another embodiment, the RVA includes
23 a virtual memory address, wherein the virtual memory address is translated into a
24 physical memory address using the RTT.

25 FIG. 6A illustrates a block diagram of a remote-translation table (RTT)
26 resident on an M chip, according to one embodiment of the present invention. M
27 chip 600 on a given node in a multi-node system includes RTT 601.

28 FIG. 6B illustrates a more detailed diagram of the RTT shown in FIG. 6A,
29 according to one embodiment on the invention. RTT 601 is indexed by a GASID
30 and high-order bits of an RVA. The GASID comprises the highest-order bits of the

1 index into RTT 601. RTT 601 is partitioned into various sections. Because the
 2 GASID comprises the highest-order bits, this embodiment shows RTT 601 being
 3 partitioned into sections corresponding to the different GASID's. In one
 4 implementation, these GASID's are associated with specific applications (or
 5 processors) operating on a MSP. RTT 601 includes section 602 for translation
 6 information corresponding to GASID₀ (at the top of RTT 601). Section 602 will not
 7 necessarily include translation information that is contiguous (i.e., in order). The
 8 information will be ordered as it is implemented by the application use for GASID₀.
 9 RTT 601 contains translation information for the entire virtual memory address
 10 space for the node on which it resides, and therefore not all of the information in
 11 section 602 is used, or contiguous (if used). Section 604 includes translation
 12 information corresponding to GASID₁, and section 606 includes translation
 13 information corresponding to GASID_A.

14 In one implementation, the index into RTT 601 includes 2 high-order bits for
 15 the GASID, and 14 high-order bits from the RVA, thereby producing a 16-bit index
 16 into RTT 601. In this implementation, there are four GASID's (from the 2 GASID
 17 bits), and therefore A is equal to 3. RTT 601 includes 64K entries (2^{16}), and each of
 18 sections 602, 604, and 606 includes 16K entries, wherein not all of the entries are
 19 necessarily applicable, or used, for the remote translation mechanism.

20 FIG. 6C illustrates a more detailed view of a portion of the RTT shown in
 21 FIG. 6B, according to one embodiment of the invention. FIG. 6C illustrates a
 22 detailed view of section 602 (corresponding to GASID₀) in RTT 601. In this
 23 embodiment, certain high-order bits of the RVA used to index into RTT 601
 24 correspond to virtual processing elements (VPE) on a given node. Translation
 25 information is ordered within section 602 according to the VPE to which it is
 26 associated. Section 608 includes information corresponding to VPE₀. Section 610
 27 includes information corresponding to VPE₁, and section 612 includes information
 28 corresponding to VPE_B.

1 In one implementation, there are 4 VPE's (in a MSP system), and therefore
 2 B is equal to 3. In this implementation, each of sections 608, 610, and 612 includes
 3 4K entries (for this portion of RTT 601).

4 FIG. 6D illustrates an even further detailed view of a portion of the RTT
 5 shown in FIG. 6C, according to one embodiment of the invention. FIG. 6D
 6 illustrates a detailed view of section 608 (corresponding to VPE₀) in RTT 601. In
 7 this embodiment, certain bits of the RVA used to index into RTT 601 correspond to
 8 segments. Information is ordered within section 608 according to the segment to
 9 which it is associated. Section 614 includes information corresponding to seg₀.
 10 Section 616 includes information corresponding to seg₁, and section 618 includes
 11 information corresponding to seg_C.

12 In some embodiments, a way to think of the remote translation mechanism is
 13 as a hardware facility that supports automatically "stitching together" the local
 14 address spaces of cooperating tasks on a series of nodes. When remote translation is
 15 enabled, the hardware routes memory access requests to the logical nodes indicated
 16 in the VNode field of the user address. FIG. 7 provides a functional view of an
 17 exemplary embodiment of this routing model.

18 Using this routing model, it is possible to devise an approach to mapping the
 19 processing-element-relative space by first building a series of independent on-node
 20 spaces, each consisting of the address spaces for their respective individual
 21 processing elements (i.e., one address space per processing element per node), and
 22 then combining them into a single contiguous space by enabling remote translation
 23 across a job or application. In some embodiments, each of the on-node address
 24 spaces are built and managed locally using source mode translation, thus allowing
 25 the use of standard memory mapping techniques and APIs for management of this
 26 space.

27 Each of these on-node address spaces are then exported (or broadcast) across
 28 the job or application space by loading the local address mappings into the RTT and
 29 entering remote translation mode. In some embodiments, it is important for all
 30 nodes to synchronize their transition to remote translation mode before resuming

1 normal execution to ensure that valid translations exist for all nodes across the job
2 or application space. In some such embodiments, failure to synchronize the
3 transition to remote translation mode results in unexpected application termination
4 because unsuccessful remote address translations are always fatal on some systems
5 and often fatal on other systems.

6 Once this address space initialization has completed and the job or
7 application is in normal operation, the operating system is responsible for
8 maintaining coherency of the remote mapping space. The operating system will
9 ensure that all valid local translations for the distributed memory job are loaded in
10 the RTT. This requires that all memory for all local processing elements within the
11 job or application is resident whenever the job is active. In addition, the operating
12 system requires that whenever the memory space for a particular processing element
13 is grown, it is grown on the processing element's local node. Doing so allows the
14 operating system to use its normal methodology for handling memory growth within
15 distributed memory jobs. After growing the processing element's address space
16 locally as requested, the RTT can be updated, making the new space available to
17 remote processing elements within the job or application.

18 FIG. 8 shows an exemplary embodiment 800 of how a distributed address
19 space is constructed according to one embodiment of the present invention. Step 1
20 of the example 800 shows the layout of the initial processing element's address
21 space. In Step 2, the initial processing element copies itself and places children, one
22 per node, across the node space of the job. Step 3 shows the continuation of the
23 example 800 by illustrating the construction of the local address and processing
24 element space on Node 0. As shown in steps 3b and 3c of this example 800, the
25 address space for the processing elements local to this node is built using standard
26 memory APIs. Step 3d then creates the remaining processing elements. The
27 example 800 keeps processing element 0's memory private during this step to allow
28 each processing element to complete its initialization in a protected address space.

29 At this point the on-node processing element and memory space
30 initialization is complete and the component node spaces are ready to combine into

1 the desired aggregate. FIG. 9 illustrates one embodiment of how this aggregation is
2 performed.

3 Once the local address spaces have been constructed on the respective nodes,
4 a representative processing element from each node group will request the operating
5 system to load the remote translation table and enable remote translation mode. FIG.
6 9 shows one embodiment of how this is performed. In step 4a the independent
7 node address spaces are shown. Note that each has been offset during its
8 initialization based on its place within the global job space and according to the
9 programming environment's distributed memory addressing model. In some
10 embodiments, at this point the processing elements on each local node have
11 synchronized with one another and are prepared to enter the global address space.
12 This is performed through an API call by each processing element to the operating
13 system that notifies the OS to enable remote translation. In some embodiments, as
14 part of this call, the OS ensures that all local memory for the job or application is
15 resident and translations are loaded in the RTT.

16 In some such embodiments utilizing synchronization, upon a successful
17 synchronization, the processing element will be executing with remote translation
18 enabled. This is illustrated at step 4b. Because each node/processing element could
19 reach this point at different times, it is prudent in some embodiments for user code
20 to synchronize across processing elements before utilizing remote memory. In some
21 embodiments, failure to do so could have undesired results including job or
22 application fatal errors or, in some other further embodiments, fatal kernel errors.

23 In some embodiments, once the job or application is in normal operation in
24 remote translation mode, the operating system will handle requests to change the
25 address space configuration on a node-local basis. For example, in some
26 embodiments, auto-grow-regions which are in the address space must be grown on
27 the node to which they were mapped. In some such embodiments this allows the
28 maintenance of normal mapping semantics by the operating system. In some
29 additional embodiments, once in remote translation mode, any attempt to modify the

1 address space of a distributed application outside scope of the local node will either
2 fail or cause the job or application to terminate.

3 One aspect of the present invention shown in FIG. 10 includes a method
4 1000 of accessing shared memory in a computer system having a plurality of nodes,
5 including a first node, wherein each node includes a processor and local memory.
6 In some embodiments, this method 1000 includes distributing 1002 an application
7 across the plurality of nodes and building 1004 an application virtual address space.
8 In some such embodiments, building 1004 an application virtual address space
9 includes building 1006 a local virtual address space for the application in each of the
10 plurality of nodes, wherein the local virtual address space translates a virtual address
11 generated by the application executing on that node to a physical address in local
12 memory for that node, and exporting 1008 the local virtual address space for each
13 node to a Remote Translation Table (RTT) associated with that node. This aspect of
14 the present invention further includes performing 1012 a memory reference to a
15 memory location in the application virtual address space, wherein performing 1012
16 a memory reference to a memory location in the application virtual address space
17 includes translating bits of the application virtual address into a node address
18 associated with the first node and translating bits of the application virtual address
19 using the RTT associated with the first node. In some further embodiments, the
20 local address space is read from a Translation Look-aside Buffer (TLB). Yet further
21 embodiments of the method 1000 include optionally performing 1010 a
22 synchronization operation that causes at least some of the plurality of nodes to wait
23 for all nodes to complete exporting their respective local virtual address spaces.

24 Another aspect of the present invention shown in FIG. 11 provides a system
25 1100 for remote address translation on a multinode system 1100 capable of
26 distributing an application, job, or process across multiple nodes 1102. In some
27 embodiments, the system 1100 includes a plurality of nodes 1102, each node 1102
28 having one or more processors 1104, a memory 1106, and a memory controller
29 1108 operatively coupled 1110 to the memory 1108 and the one or more
30 processors 1104. In some such embodiments, the memory controller 1108 includes a

1 Remote Translation Table (RTT) 1112, wherein the RTT 1112 translates a virtual
2 address received as part of a memory request received from another node 1102 into
3 a memory request with physical addresses into the memory 1106 on the node 1102
4 associated with the RTT 1112. Further within some of these embodiments, the RTT
5 1112 is initialized upon the start of an process associated with an application or
6 process by building a virtual to physical address translations for local virtual address
7 space in the node 1102 corresponding to the application, and exporting the virtual to
8 physical address translations for the local virtual address space from the node 1102
9 to the RTT 1112 associated with that node 1102. In some embodiments, each of the
10 plurality of nodes 1102 executes a synchronization operation that causes at least
11 some of the plurality of nodes 1102 to wait for all of the plurality of nodes 1102 to
12 complete exporting the virtual to physical address translations to their respective
13 RTT's 1112.

14 Yet another aspect of the present invention provides a device-readable
15 medium having instructions thereon that cause a properly programmed device to
16 perform a method of accessing shared memory in the device. In some
17 embodiments, the instructions, when executed on a properly programmed
18 information-processing device having a plurality of nodes, including a first node,
19 each node having one or more processors, a memory, and a memory controller and
20 coupled to the memory and the one or more processors, cause the information-
21 processing device to distribute an application across the plurality of nodes and build
22 an application virtual address space. In some such embodiments, building an
23 application virtual address space includes building a local virtual address space for
24 the application in each of the plurality of nodes, wherein the local virtual address
25 space translates a virtual address generated by the application executing on that
26 node to a physical address in local memory for that node, and exporting the local
27 virtual address space for each node to a Remote Translation Table (RTT) associated
28 with that node. In some embodiments, the instructions, when executed further
29 include performing a memory reference to a memory location in the application
30 virtual address space, wherein performing a memory reference to a memory location

1 in the application virtual address space includes translating bits of the application
2 virtual address into a node address associated with the first node and translating bits
3 of the application virtual address using the RTT associated with the first node. In
4 some embodiments, building a local virtual address space further includes
5 performing a synchronization operation that causes at least some of the plurality of
6 nodes to wait for all nodes complete exporting their respective address space. In
7 some further embodiments, the local address space is read from a Translation Look-
8 aside Buffer (TLB).

9 Still another aspect of the present invention provides a multinode system for
10 implementing remote address translation. Some embodiments of the multinode
11 system include a plurality of nodes, including a first node. In some such
12 embodiments, each of the plurality of nodes includes one or more processors, a
13 memory, and a memory controller operatively coupled to the memory and the one or
14 more processors. These embodiments include a means for distributing an
15 application across the plurality of nodes and a means for building an application
16 virtual address space. In various embodiments, the means for building an
17 application virtual address space includes a means for building a local virtual
18 address space for the application in each of the plurality of nodes, wherein the local
19 virtual address space translates a virtual address generated by the application
20 executing on that node to a physical address in local memory for that node and a
21 means for exporting the local virtual address space for each node to a Remote
22 Translation Table (RTT) associated with that node. Some further embodiments
23 include a means for performing a memory reference to a memory location in the
24 application virtual address space, wherein performing a memory reference to a
25 memory location in the application virtual address space includes a means for
26 translating bits of the application virtual address into a node address associated with
27 the first node, and a means for translating bits of the application virtual address
28 using the RTT associated with the first node. In some such embodiments, building
29 an application virtual address space further includes a means for performing a
30 synchronization operation that causes at least some of the plurality of nodes to wait

1 for all nodes to complete exporting their respective local virtual address spaces.

2

3 As described herein, the various embodiments of the present invention
4 provide a number of advantages. For example, an RTT provides a scalable address
5 translation mechanism, and is designed to avoid translation faults in large systems
6 (unlike a regular TLB design). The RTT supports full mapping of all the memory in
7 a machine (unlike various networking cards) to allow full load/store access to all the
8 memory in the system. Such a system allows each node to independently manage
9 its own virtual-to-physical memory mapping. Such a system also removes the need
10 to implement conventional TLB "shutdown." Conventional TLB "shutdown"
11 occurs when a node changes a local virtual-to-physical page mapping, and has to
12 invalidate all of the TLB entries throughout the system that contain that mapping.
13 The use of an RTT that supports full mapping removes the need to implement such
14 an approach. These and other advantages are provided for by various embodiments
15 of the present invention.

16 It is understood that the above description is intended to be illustrative, and
17 not restrictive. Many other embodiments will be apparent to those of skill in the art
18 upon reviewing the above description. The scope of the invention should, therefore,
19 be determined with reference to the appended claims, along with the full scope of
20 equivalents to which such claims are entitled.